

# SDS 383D Chapter 1: Preliminaries

Xizewen Han

March 11, 2019

## Bayesian Inference in Simple Conjugate Families

(A)

*Sol:* For  $i = 1, \dots, N$  such that  $x_i \sim \text{Bernoulli}(w)$ , we have

$$p(x_i|w) = w^{x_i} (1-w)^{1-x_i}$$

Therefore with prior  $w \sim \text{Beta}(a, b)$ , we have the posterior distribution

$$\begin{aligned} p(w|x_{1:N}) &\propto p(x_{1:N}|w) \cdot p(w) \\ &= \prod_{i=1}^N p(x_i|w) \cdot p(w) \\ &= w^{\sum x_i} (1-w)^{N-\sum x_i} \cdot \frac{1}{B(a, b)} w^{a-1} (1-w)^{b-1} \\ &\propto w^{\sum x_i + a - 1} (1-w)^{N - \sum x_i + b - 1} \end{aligned}$$

Note that the last expression is the kernel of  $\text{Beta}(\sum_{i=1}^N x_i + a, N - \sum_{i=1}^N x_i + b)$  distribution, which by adding the normalizing term we'd have this as the posterior distribution. ■

(B)

*Sol:* For such transformation, which is one-to-one and onto, we have the inverse as

$$\begin{cases} x_1 = y_1 y_2, \\ x_2 = y_2 - y_1 y_2 \end{cases}$$

Therefore, we have the Jacobian as

$$J = \begin{bmatrix} \frac{\partial x_1}{\partial y_1} & \frac{\partial x_1}{\partial y_2} \\ \frac{\partial x_2}{\partial y_1} & \frac{\partial x_2}{\partial y_2} \end{bmatrix} = \begin{bmatrix} y_2 & y_1 \\ -y_2 & 1 - y_1 \end{bmatrix} = y_2 \geq 0$$

thus the joint distribution is

$$\begin{aligned}
f_{Y_1, Y_2}(y_1, y_2) &= f_{X_1, X_2}(x_1(y_1, y_2), x_2(y_1, y_2))|J| \\
&= f_{X_1}(x_1(y_1, y_2))f_{X_2}(x_2(y_1, y_2))|J| \\
&= \frac{1}{\Gamma(a_1)}(y_1 y_2)^{a_1-1} e^{-y_1 y_2} \cdot \frac{1}{\Gamma(a_2)}(y_2 - y_1 y_2)^{a_2-1} e^{-y_2 + y_1 y_2} \cdot y_2 \\
&= \underbrace{\frac{1}{\Gamma(a_1)\Gamma(a_2)} y_1^{a_1-1} (1-y_1)^{a_2-1}}_{g(y_1), \text{ kernel of Beta}(a_1, a_2)} h(y_2), \text{ kernel of Gamma}(a_1+a_2, 1)
\end{aligned}$$

Since the joint distribution of  $Y_1$  and  $Y_2$  can be factored into two functions,  $g(y_1)$  and  $h(y_2)$ , by Lemma 4.2.7 of Casella and Berger's Statistical Inference, as well as the form of each function noted above in the last expression, we conclude that  $Y_1$  and  $Y_2$  are independent random variables such that

$$Y_1 \sim \text{Beta}(a_1, a_2)$$

and

$$Y_2 \sim \text{Gamma}(a_1 + a_2, 1)$$

Based on such conclusion, we propose the following method to simulate a  $\text{Beta}(a_1, a_2)$  random variable:

- 1st step: simulate from two Gamma random variables:  $X_1 \sim \text{Gamma}(a_1, 1)$  and  $X_2 \sim \text{Gamma}(a_2, 1)$
- 2nd step: calculate  $Y_1 = \frac{X_1}{X_1 + X_2}$

Such  $Y_1 \sim \text{Beta}(a_1, a_2)$ . ■

### (C)

*Sol:* With  $X_i \sim \mathcal{N}(\theta, \sigma^2)$  for  $i = 1, \dots, N$ , in which  $\theta$  is unknown and  $\sigma^2$  is known, and  $\theta \sim \mathcal{N}(m, v)$ , we have the posterior distribution

$$\begin{aligned}
p(\theta|x_{1:N}) &\propto p(\theta) \cdot p(x_{1:N}|\theta) \\
&= \frac{1}{\sqrt{2\pi v}} e^{-\frac{(\theta-m)^2}{2v}} \cdot \prod_{i=1}^N \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(x_i-\theta)^2}{2\sigma^2}} \\
&\propto \exp\left(-\frac{(nv + \sigma^2)\theta^2 - 2(m\sigma^2 + v(\sum x_i))}{2\sigma^2 v}\right)
\end{aligned}$$

By noting the kernel, we recognize the posterior distribution as

$$\mathcal{N}\left(\frac{m\sigma^2 + v(\sum_{i=1}^N x_i)}{\sigma^2 + nv}, \frac{\sigma^2 v}{\sigma^2 + nv}\right)$$

The variance term can be re-written as

$$\frac{\sigma^2 v}{\sigma^2 + nv} = \frac{1}{\frac{1}{v} + \frac{n}{\sigma^2}}$$

Note the denominator is the precision, and it's the sum of data precision and prior prevision – a good characteristic of posterior of normal-normal model. ■

(D)

*Sol:* With  $X_i \sim \mathcal{N}(\theta, 1/w)$  for  $i = 1, \dots, N$ , in which  $\theta$  is known and precision  $w$  is unknown, and  $w \sim \text{Gamma}(a, b)$ , we have the posterior distribution

$$\begin{aligned} p(w|x_{1:N}) &\propto p(w) \cdot \prod_{i=1}^N p(x_i|\theta, w) \\ &\propto w^{a-1} e^{-\frac{1}{b}w} w^{\frac{N}{2}} e^{-\frac{\sum(x_i-\theta)^2}{2}w} \\ &= w^{\frac{N}{2}+a-1} e^{-\left(\frac{1}{b} + \frac{\sum(x_i-\theta)^2}{2}\right)w} \end{aligned}$$

By recognizing the kernel, we have the posterior distribution

$$w \sim \text{Gamma}\left(\frac{N}{2} + a, \frac{1}{b} + \frac{\sum(x_i - \theta)^2}{2}\right)$$

Since  $w$  is a re-parameterization of  $\sigma^2$ , we use change of variables directly on posterior pdf to derive the posterior of  $\sigma^2$ . We have the transformation  $g(w) = \frac{1}{\sigma^2}$ , which is one-to-one and onto, and

$$\frac{d}{d\sigma^2} g^{-1}(\sigma^2) = \frac{d}{d\sigma^2} (\sigma^2)^{-1} = -(\sigma^2)^{-2}$$

Therefore, for a general case  $w \sim \text{Gamma}(a, b)$ , we have

$$p(\sigma^2) = p(g^{-1}(\sigma^2)) \left| \frac{d}{d\sigma^2} g^{-1}(\sigma^2) \right| = \frac{b^a}{\Gamma(a)} (\sigma^2)^{-(a-1)} e^{-\frac{b}{\sigma^2}} (\sigma^2)^{-2} = \frac{b^a}{\Gamma(a)} (\sigma^2)^{-a-1} e^{-\frac{b}{\sigma^2}}$$

i.e.  $\sigma^2 \sim IG(a, b)$ . Thus for our specific case, the posterior distribution of  $\sigma^2$  is

$$\sigma^2 \sim IG\left(\frac{N}{2} + a, \frac{1}{b} + \frac{\sum(x_i - \theta)^2}{2}\right)$$

■

(E)

*Sol:* With known idiosyncratic variances for normal likelihood, the posterior density for unknown common mean  $\theta$  is

$$\begin{aligned} p(\theta|x_{1:N}) &\propto \mathbb{P}(\theta) \cdot \prod_{i=1}^N p(x_i|\theta) \\ &= \frac{1}{\sqrt{2\pi v}} e^{-\frac{(\theta-m)^2}{2v}} \frac{1}{\sqrt{2\pi\sigma_i^2}} e^{-\frac{1}{2} \sum \frac{(x_i-\theta)^2}{\sigma_i^2}} \\ &\propto e^{-\frac{1}{2} \left( \left( \frac{1}{v} + \sum \frac{1}{\sigma_i^2} \right) \theta^2 - 2 \left( \frac{m}{v} + \sum \frac{x_i}{\sigma_i^2} \right) \theta \right)} \end{aligned}$$

We can thus directly write out the posterior distribution for  $\theta$  as

$$\mathcal{N}\left(\frac{\frac{1}{v}m + \sum_{i=1}^N \frac{1}{\sigma_i^2}x_i}{\frac{1}{v} + \sum_{i=1}^N \frac{1}{\sigma_i^2}}, \frac{1}{v} + \sum_{i=1}^N \frac{1}{\sigma_i^2}\right)$$

notice that the mean term is expressed in the form as a weighted average of the observations  $x_i$ 's and the prior mean  $m$ .

A later note: the weighted average part is more obvious if we work with precision instead of variance: denote prior precision as  $w = \frac{1}{v}$  and idiosyncratic precision for each observation as  $w_i = \frac{1}{\sigma_i^2}$ , we have the posterior distribution in the form of

$$\mathcal{N}\left(\frac{w \cdot m + \sum_{i=1}^N w_i \cdot x_i}{w + \sum_{i=1}^N w_i}, w + \sum_{i=1}^N w_i\right)$$

■

(F)

*Sol:* Using the parameterization from the question sheet, we compute the marginal distribution of  $x$  as

$$\begin{aligned} p(x) &= \int_w p(x|w)p(w)dw \\ &= \int_w \left(\frac{w}{2\pi}\right)^{1/2} e^{-\frac{w}{2}(x-m)^2} \frac{(b/2)^{a/2}}{\Gamma(a/2)} w^{a/2-1} e^{-\frac{b}{2}w} dw \\ &= \frac{(b/2)^{a/2}}{\sqrt{2\pi}\Gamma(a/2)} \cdot \int_0^\infty \underbrace{w^{\frac{a+1}{2}-1} e^{-\left(\frac{b}{2} + \frac{(x-m)^2}{2}\right)w}}_{\text{kernel of } \text{Gamma}(\frac{a+1}{2}, \frac{b+(x-m)^2}{2})} dw \\ &= \frac{(b/2)^{a/2}}{\sqrt{2\pi}\Gamma(a/2)} \frac{\Gamma(\frac{a+1}{2})}{\left(\frac{b+(x-m)^2}{2}\right)^{\frac{a+1}{2}}} \cdot 1 \\ &= \frac{\Gamma(\frac{a+1}{2})}{\Gamma(\frac{a}{2})} \frac{1}{\sqrt{\pi b}} \left(1 + \frac{(x-m)^2}{b}\right)^{-\frac{a+1}{2}} \\ &= \frac{\Gamma(\frac{a+1}{2})}{\Gamma(\frac{a}{2})} \frac{1}{\sqrt{a\pi(b/a)^{1/2}}} \left(\frac{a + \left(\frac{x-m}{(b/a)^{1/2}}\right)^2}{a}\right)^{-\frac{a+1}{2}} \end{aligned}$$

Note that the last expression is exactly the pdf of  $t$  location-scale distribution with  $d = a$  degrees of freedom, center  $m$ , and scale parameter  $(b/a)^{1/2}$ .

Later note: this question demonstrates a useful conclusion – if we want to sample from  $t$ -distribution, which has heavy tails, we don't have to sample directly; instead, we first sample precision from Gamma distribution, then sample observation from normal distribution given that precision. The resulting sample is a random variable with  $t$ -distribution.

■

## The Multivariate Normal Distribution

(A)

*Sol:* For the two sub-problems, we do them by using the definition given by the question and directly applying matrix operation rules:

(1)

$$\begin{aligned}
\text{cov}(\mathbf{x}) &= \mathbb{E}((\mathbf{x} - \boldsymbol{\mu})(\mathbf{x} - \boldsymbol{\mu})^T) \\
&= \mathbb{E}((\mathbf{x} - \boldsymbol{\mu})(\mathbf{x}^T - \boldsymbol{\mu}^T)) \\
&= \mathbb{E}(\mathbf{x}\mathbf{x}^T - \mathbf{x}\boldsymbol{\mu}^T - \boldsymbol{\mu}\mathbf{x}^T + \boldsymbol{\mu}\boldsymbol{\mu}^T) \\
&= \mathbb{E}(\mathbf{x}\mathbf{x}^T) - \mathbb{E}(\mathbf{x})\boldsymbol{\mu}^T - \boldsymbol{\mu}\mathbb{E}(\mathbf{x}^T) + \boldsymbol{\mu}\boldsymbol{\mu}^T \\
&= \mathbb{E}(\mathbf{x}\mathbf{x}^T) - \boldsymbol{\mu}\boldsymbol{\mu}^T - \boldsymbol{\mu}\boldsymbol{\mu}^T + \boldsymbol{\mu}\boldsymbol{\mu}^T \text{ (using the fact that } \mathbb{E}(\mathbf{x}) = \boldsymbol{\mu}) \\
&= \mathbb{E}(\mathbf{x}\mathbf{x}^T) - \boldsymbol{\mu}\boldsymbol{\mu}^T
\end{aligned}$$

The last expression is what's asked to show. ■

(2)

*Sol:* The mean of  $\mathbf{A}\mathbf{x} + \mathbf{b}$  is:

$$\mathbb{E}(\mathbf{A}\mathbf{x} + \mathbf{b}) = \mathbf{A}\mathbb{E}(\mathbf{x}) + \mathbb{E}(\mathbf{b}) = \mathbf{A}\boldsymbol{\mu} + \mathbf{b}$$

Therefore, we have

$$\begin{aligned}
\text{cov}(\mathbf{A}\mathbf{x} + \mathbf{b}) &= \mathbb{E}\left(\left((\mathbf{A}\mathbf{x} + \mathbf{b}) - (\mathbf{A}\boldsymbol{\mu} + \mathbf{b})\right)\left((\mathbf{A}\mathbf{x} + \mathbf{b}) - (\mathbf{A}\boldsymbol{\mu} + \mathbf{b})\right)^T\right) \\
&= \mathbb{E}\left(\left(\mathbf{A}(\mathbf{x} - \boldsymbol{\mu})\right)\left(\mathbf{A}(\mathbf{x} - \boldsymbol{\mu})\right)^T\right) \\
&= \mathbb{E}\left(\mathbf{A}(\mathbf{x} - \boldsymbol{\mu})(\mathbf{x} - \boldsymbol{\mu})^T \mathbf{A}^T\right) \\
&= \mathbf{A}\mathbb{E}\left((\mathbf{x} - \boldsymbol{\mu})(\mathbf{x} - \boldsymbol{\mu})^T\right)\mathbf{A}^T \\
&= \mathbf{A}\text{cov}(\mathbf{x})\mathbf{A}^T \text{ (the middle term is derived by the definition of covariance given by the question)}
\end{aligned}$$

The last expression is what's asked to show. ■

(B)

*Sol:* Since for all  $z_i$ 's such that  $i = 1, \dots, p$ , they are independent, the pdf of  $\mathbf{z} = (z_1, \dots, z_p)^T$ , being the joint distribution of  $z_i$ 's, is the multiplication of their individual densities:

$$\begin{aligned}
f_{\mathbf{Z}}(\mathbf{z}) &= \prod_{i=1}^p f(z_i) \\
&= \prod_{i=1}^p \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{z_i^2}{2}\right) \\
&= \frac{1}{(2\pi)^{p/2}} \cdot \exp\left(-\frac{1}{2} \sum z_i^2\right) \\
&= \frac{1}{(2\pi)^{p/2}} \cdot \exp\left(-\frac{1}{2} \mathbf{z}^T \mathbf{z}\right)
\end{aligned}$$

The mgf of  $\mathbf{z}$  is given by

$$\begin{aligned}
\mathbb{E}(e^{\mathbf{t}^T \mathbf{z}}) &= \int_{z_1} \cdots \int_{z_p} e^{\mathbf{t}^T \mathbf{z}} f_{\mathbf{Z}}(\mathbf{z}) dz_1 \dots dz_p \\
&= \int_{z_1} \cdots \int_{z_p} e^{\sum t_i z_i} f_{\mathbf{Z}}(\mathbf{z}) dz_1 \dots dz_p \\
&= \int_{z_1} \cdots \int_{z_p} \prod_{i=1}^p e^{t_i z_i} \prod f_{Z_i}(z_i) dz_1 \dots dz_p \\
&= \prod_{i=1}^p \int_{z_i} e^{t_i z_i} f_{Z_i}(z_i) dz_i \text{ (using the fact that } z_i \text{'s are independent)} \\
&= \prod_{i=1}^p \exp(m t_i + v t_i^2 / 2) \text{ (given by question)} \\
&= \prod_{i=1}^p \exp(0 \cdot t_i + 1 \cdot t_i^2 / 2) \text{ (given by question)} \\
&= \prod_{i=1}^p \exp(t_i^2 / 2) \\
&= \exp\left(\frac{t^2}{2}\right) \\
&= \exp\left(\frac{1}{2} \mathbf{t}^T \mathbf{t}\right)
\end{aligned}$$

■

## (C)

*Sol:* To show that if  $\mathbf{x} \sim \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ , its mgf is of the given form, we use the definition given in this problem: for a multivariate normal variable  $\mathbf{x}$ ,  $z = \mathbf{a}^T \mathbf{x}$  is univariate normal.

It's expectation is

$$\mathbb{E}(z) = \mathbb{E}(\mathbf{a}^T \mathbf{x}) = \mathbf{a}^T \mathbb{E}(\mathbf{x}) = \mathbf{a}^T \boldsymbol{\mu}$$

and

$$var(z) = cov(\mathbf{a}^T \mathbf{x}) = \mathbf{a}^T cov(\mathbf{a}^T \mathbf{x}) \mathbf{a} = \mathbf{a}^T \boldsymbol{\Sigma} \mathbf{a}$$

We know the mgf of a univariate normal variable as

$$M_z(t) = (\mathbf{a}^T \boldsymbol{\mu}) t + (\mathbf{a}^T \boldsymbol{\Sigma} \mathbf{a}) \frac{t^2}{2} = (t \mathbf{a})^T \boldsymbol{\mu} + \frac{(t \mathbf{a})^T \boldsymbol{\Sigma} (t \mathbf{a})}{2}$$

meanwhile

$$M_z(t) = \mathbb{E}(e^{zt}) = \mathbb{E}(e^{(\mathbf{a}^T \mathbf{x})t}) = \mathbb{E}(e^{(t \mathbf{a})^T \mathbf{x}}) = M_{\mathbf{x}}(t \mathbf{a})$$

therefore, denoting  $\mathbf{t} = t \mathbf{a}$ , we have the mgf of multivariate normal variable  $\mathbf{x}$  as

$$M_{\mathbf{x}}(\mathbf{t}) = \mathbf{t}^T \boldsymbol{\mu} + \frac{\mathbf{t}^T \boldsymbol{\Sigma} \mathbf{t}}{2}$$

For the other direction, i.e. with the given form of mgf, such variable is multivariate normal, we again use the definition given in this particular question s.t. every linear combination of a multivariate normal random variable is univariate normal:

for all vectors  $\mathbf{a}$  such that  $z = \mathbf{a}^T \mathbf{x}$ , we have its moment generating function as

$$\begin{aligned} M_z(t) &= M_{\mathbf{a}^T \mathbf{x}}(t) \\ &= \mathbb{E}(e^{t \cdot (\mathbf{a}^T \mathbf{x})}) \\ &= \mathbb{E}(e^{(t\mathbf{a})^T \cdot \mathbf{x}}) \\ &= M_{\mathbf{x}}(t\mathbf{a}) \\ &= \exp((t\mathbf{a})^T \boldsymbol{\mu} + \frac{1}{2}(t\mathbf{a})^T \boldsymbol{\Sigma}(t\mathbf{a})) \\ &= \exp(t(\mathbf{a}^T \boldsymbol{\mu}) + \frac{1}{2}t^2(\mathbf{a}^T \boldsymbol{\Sigma} \mathbf{a})) \\ &= \exp(m \cdot t + v \cdot t^2/2) \end{aligned}$$

the last expression is exactly the form of the mgf of univariate normal distribution  $\mathcal{N}(m, v)$ , in which  $m = \mathbf{a}^T \boldsymbol{\mu}$ ,  $v = \mathbf{a}^T \boldsymbol{\Sigma} \mathbf{a}$ , i.e.  $z$  has a  $\mathcal{N}(m, v)$  distribution. Therefore, according to the definition given by the question,  $\mathbf{x}$  is multivariate normal.  $\blacksquare$

## (D)

*Sol:* The question defines

$$\mathbf{x} = \mathbf{L}\mathbf{z} + \boldsymbol{\mu}$$

in which

$$\mathbf{z} = (z_1, \dots, z_p)^T$$

s.t. each  $z_i$  for  $i = 1, \dots, p$  are i.i.d.  $\mathcal{N}(0, 1)$  distribution.

I want to solve this problem with two methods, with the intention of going through different parts of mathematical statistics knowledge as a review:

### Method 1

Since  $z_i$ 's are mutually independent, by Corollary 4.6.10 from Casella and Berger's *Statistical Inference* each of their linear combinations is univariate normal, thus  $\mathbf{z}$  is a multivariate normal random variable. Its pdf can be derived by

$$p(\mathbf{z}) = \prod_{i=1}^p \frac{1}{\sqrt{2\pi}} e^{-z_i^2/2} = \frac{1}{(2\pi)^{p/2}} e^{-\mathbf{z}^T \mathbf{z}/2}$$

which is the form of a standard multivariate normal variable's pdf – we just proved the claim/definition given from (B). We directly write out our result from (B): the mgf of  $\mathbf{z}$  is

$$M_z(\mathbf{t}) = \exp\left(\frac{1}{2}\mathbf{t}^T \mathbf{t}\right)$$

Now we show  $\mathbf{x}$  is a multivariate normal variable: using change of variable, we have

$$\mathbf{x} = g(\mathbf{z}) = \mathbf{L}\mathbf{z} + \boldsymbol{\mu}$$

since  $\mathbf{L}$  is full rank, it's invertible, thus we have

$$\mathbf{z} = g^{-1}(\mathbf{x}) = \mathbf{L}^{-1}(\mathbf{x} - \boldsymbol{\mu})$$

Therefore we have the Jacobian

$$J = \left| \frac{\partial \mathbf{r}(\mathbf{x})}{\partial \mathbf{x}} \right| = \left| \frac{\partial \mathbf{g}^{-1}(\mathbf{x})}{\partial \mathbf{x}} \right| = |\mathbf{L}^{-1}| = |\mathbf{L}|^{-1}$$

The pdf of  $\mathbf{x}$  is thus

$$\begin{aligned} p(\mathbf{x}) &= \frac{1}{(2\pi)^{p/2}} e^{-(\mathbf{L}^{-1}(\mathbf{x} - \boldsymbol{\mu}))^T (\mathbf{L}^{-1}(\mathbf{x} - \boldsymbol{\mu})) / 2} \cdot |J| \\ &= (2\pi)^{-p/2} |\mathbf{L}|^{-1} \exp\left(-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^T (\mathbf{L}^{-1})^T (\mathbf{L}^{-1})(\mathbf{x} - \boldsymbol{\mu})\right) \\ &= (2\pi)^{-p/2} |\mathbf{L}|^{-1} \exp\left(-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^T (\mathbf{L}^T)^{-1} (\mathbf{L}^{-1})(\mathbf{x} - \boldsymbol{\mu})\right) \\ &= (2\pi)^{-p/2} |\mathbf{L}|^{-1} \exp\left(-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^T (\mathbf{L}\mathbf{L}^T)^{-1}(\mathbf{x} - \boldsymbol{\mu})\right) \end{aligned}$$

Since  $\mathbf{L}$  is full rank,  $\mathbf{L}\mathbf{L}^T = \mathbf{L}\mathbf{I}_{p \times p}\mathbf{L}^T$  is positive definite, in which  $\mathbf{I}_{p \times p}$  is identity matrix (thus positive definite). Note that  $p(\mathbf{x})$  is exactly the form of the pdf of  $\mathcal{MVN}(\boldsymbol{\mu}, \mathbf{L}\mathbf{L}^T)$  distribution.

## Method 2

I also want to use (a simple manipulation on) mgf to solve this question:

Since we've already derived the mgf of  $\mathbf{z}$  as a standard multivariate normal random variable, i.e.

$$M_z(\mathbf{t}) = \exp\left(\frac{1}{2}\mathbf{t}^T \mathbf{t}\right)$$

we can directly write out the mgf of  $\mathbf{x}$  as

$$\begin{aligned} M_{\mathbf{x}}(\mathbf{t}) &= \mathbb{E}(e^{\mathbf{t}^T \mathbf{x}}) \\ &= \mathbb{E}(e^{\mathbf{t}^T (\mathbf{L}\mathbf{z} + \boldsymbol{\mu})}) \\ &= \exp(\mathbf{t}^T \boldsymbol{\mu}) \mathbb{E}((\mathbf{L}^T \mathbf{t})^T \mathbf{z}) \\ &= \exp(\mathbf{t}^T \boldsymbol{\mu}) M_z(\mathbf{L}^T \mathbf{t}) \\ &= \exp(\mathbf{t}^T \boldsymbol{\mu}) \exp\left(\frac{1}{2}(\mathbf{L}^T \mathbf{t})^T (\mathbf{L}^T \mathbf{t})\right) \\ &= \exp\left(\mathbf{t}^T \boldsymbol{\mu} + \frac{1}{2}\mathbf{t}^T \mathbf{L}\mathbf{L}^T \mathbf{t}\right) \end{aligned}$$

thus by our conclusion from (C), a random variable with this form of mgf has a  $\mathcal{MVN}(\boldsymbol{\mu}, \mathbf{L}\mathbf{L}^T)$  distribution.

## A Later Note

We found out that the restriction such that  $\mathbf{L}$  shall be a squared matrix can be relaxed – it can be a non-squared matrix, as long as it's full row rank (not that it's not full column rank!). We'd use this finding as one of our methods in solving Question (A) from **Conditionals and Marginals**.

To compute the expected value and covariance matrix of  $\mathbf{x}$ , we take the first and second derivative of mgf at  $t = 0$ , respectively:

$$\begin{aligned}
\mathbb{E}(\mathbf{x}) &= \frac{\partial M_{\mathbf{x}}(\mathbf{t})}{\partial \mathbf{t}} \Big|_{t=0} \\
&= \frac{\partial}{\partial \mathbf{t}} \int_{\mathbf{x}} \exp(\mathbf{t}^T \boldsymbol{\mu} + \frac{1}{2} \mathbf{t}^T \mathbf{L} \mathbf{L}^T \mathbf{t}) f(\mathbf{x}) d\mathbf{x} \Big|_{t=0} \\
&= \int_{\mathbf{x}} \frac{\partial}{\partial \mathbf{t}} \exp(\mathbf{t}^T \boldsymbol{\mu} + \frac{1}{2} \mathbf{t}^T \mathbf{L} \mathbf{L}^T \mathbf{t}) f(\mathbf{x}) d\mathbf{x} \Big|_{t=0} \\
&= \int_{\mathbf{x}} \frac{\partial}{\partial \mathbf{t}} \exp(\mathbf{t}^T \boldsymbol{\mu} + \frac{1}{2} \mathbf{t}^T \mathbf{L} \mathbf{L}^T \mathbf{t}) f(\mathbf{x}) d\mathbf{x} \Big|_{t=0} \\
&= \int_{\mathbf{x}} \exp(\mathbf{t}^T \boldsymbol{\mu} + \frac{1}{2} \mathbf{t}^T \mathbf{L} \mathbf{L}^T \mathbf{t}) (\boldsymbol{\mu} + \mathbf{L} \mathbf{L}^T \mathbf{t}) f(\mathbf{x}) d\mathbf{x} \Big|_{t=0} \\
&= \int_{\mathbf{x}} \exp(\mathbf{0}) (\boldsymbol{\mu}) f(\mathbf{x}) d\mathbf{x} \\
&= \boldsymbol{\mu} \int_{\mathbf{x}} f(\mathbf{x}) d\mathbf{x} \\
&= \boldsymbol{\mu}
\end{aligned}$$

and

$$\begin{aligned}
\mathbb{E}(\mathbf{x}^T \mathbf{x}) &= \frac{\partial^2 M_{\mathbf{x}}(\mathbf{t})}{\partial \mathbf{t}^T \partial \mathbf{t}} \Big|_{t=0} \\
&= \frac{\partial}{\partial \mathbf{t}^T} \frac{\partial M_{\mathbf{x}}(\mathbf{t})}{\partial \mathbf{t}} \Big|_{t=0} \\
&= \frac{\partial}{\partial \mathbf{t}^T} \int_{\mathbf{x}} \exp(\mathbf{t}^T \boldsymbol{\mu} + \frac{1}{2} \mathbf{t}^T \mathbf{L} \mathbf{L}^T \mathbf{t}) (\boldsymbol{\mu} + \mathbf{L} \mathbf{L}^T \mathbf{t}) f(\mathbf{x}) d\mathbf{x} \Big|_{t=0} \\
&= \int_{\mathbf{x}} \frac{\partial}{\partial \mathbf{t}^T} \exp(\mathbf{t}^T \boldsymbol{\mu} + \frac{1}{2} \mathbf{t}^T \mathbf{L} \mathbf{L}^T \mathbf{t}) (\boldsymbol{\mu} + \mathbf{L} \mathbf{L}^T \mathbf{t}) f(\mathbf{x}) d\mathbf{x} \Big|_{t=0} \\
&= \int_{\mathbf{x}} \left( \exp(\mathbf{t}^T \boldsymbol{\mu} + \frac{1}{2} \mathbf{t}^T \mathbf{L} \mathbf{L}^T \mathbf{t}) (\boldsymbol{\mu} + \mathbf{L} \mathbf{L}^T \mathbf{t})^T (\boldsymbol{\mu} + \mathbf{L} \mathbf{L}^T \mathbf{t}) + \exp(\mathbf{t}^T \boldsymbol{\mu} + \frac{1}{2} \mathbf{t}^T \mathbf{L} \mathbf{L}^T \mathbf{t}) (\mathbf{L} \mathbf{L}^T) \right) f(\mathbf{x}) d\mathbf{x} \Big|_{t=0} \\
&= \int_{\mathbf{x}} \left( \exp(\mathbf{0}) (\boldsymbol{\mu})^T (\boldsymbol{\mu}) + \exp(\mathbf{0}) (\mathbf{L} \mathbf{L}^T) \right) f(\mathbf{x}) d\mathbf{x} \\
&= (\boldsymbol{\mu}^T \boldsymbol{\mu} + \mathbf{L} \mathbf{L}^T) \int_{\mathbf{x}} f(\mathbf{x}) d\mathbf{x} \\
&= \boldsymbol{\mu}^T \boldsymbol{\mu} + \mathbf{L} \mathbf{L}^T
\end{aligned}$$

therefore

$$cov(\mathbf{x}) = \mathbb{E}(\mathbf{x}^T \mathbf{x}) - \mathbb{E}(\mathbf{x})^T \mathbb{E}(\mathbf{x}) = \boldsymbol{\mu}^T \boldsymbol{\mu} + \mathbf{L} \mathbf{L}^T - \boldsymbol{\mu}^T \boldsymbol{\mu} = \mathbf{L} \mathbf{L}^T$$

■

(E)

*Sol:* For the "only if" direction, i.e. for multivariate normal variable  $\mathbf{x}$ , we could argue in the following manner: according to Casella and Berger's arguments under Theorem 2.3.12 in *Statistical Inference* about the uniqueness of Laplace transforms, each mgf can only have one pdf, i.e. one distinct distribution being

mapped to it. Therefore, the logic to solve this problem is that if we can find an affine transformation of independent univariate normal variables such that they have the same mgf as our desired multivariate normal variable's mgf, we can conclude that the transformation and the multivariate variable are from the same distribution, i.e. the multivariate normal variable can be expressed as an affine transformation of a standard multivariate normal variable.

We had from Question (C) such that for a  $\mathcal{MVN}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$  random variable  $\mathbf{x}$ , it has mgf

$$M_{\mathbf{x}}(\mathbf{t}) = \exp \left( \mathbf{t}^T \boldsymbol{\mu} + \frac{\mathbf{t}^T \boldsymbol{\Sigma} \mathbf{t}}{2} \right)$$

Meanwhile, we've found in (D) s.t. an affine transformation  $\mathbf{L}\mathbf{z} + \boldsymbol{\mu}$  to standard multivariate normal variable  $\mathbf{z}$  has mgf

$$M_{\mathbf{L}\mathbf{z} + \boldsymbol{\mu}}(\mathbf{t}) = \exp \left( \mathbf{t}^T \boldsymbol{\mu} + \frac{\mathbf{t}^T \mathbf{L} \mathbf{L}^T \mathbf{t}}{2} \right)$$

Therefore, by simply observing the forms of these two mgf's we can find, by setting  $\mathbf{L}$  s.t.  $\mathbf{L}\mathbf{L}^T = \boldsymbol{\Sigma}$ , we'd have the same form of mgf. We can achieve this by performing Cholesky decomposition on  $\boldsymbol{\Sigma}$ : since  $\boldsymbol{\Sigma}$  is positive definite, we can find  $\mathbf{L}$  s.t.  $\boldsymbol{\Sigma} = \mathbf{L}\mathbf{L}^T$ .

We've thus found such affine transformation that  $\mathbf{x} = \mathbf{L}\mathbf{z} + \boldsymbol{\mu}$ , in which  $\boldsymbol{\mu}$  is the mean of  $\mathbf{x}$  and  $\mathbf{L}$  is the component of Cholesky decomposition of covariance matrix of  $\mathbf{x}$ .

With this conclusion, I propose the following algorithm to simulate a  $\mathcal{MVN}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$  variable  $\mathbf{x}$ :

- (1) simulate  $\mathbf{z}$  from  $\mathcal{MVN}(\mathbf{0}, \mathbf{I})$  distribution
- (2) perform Cholesky decomposition on covariance matrix  $\boldsymbol{\Sigma}$  of desired multivariate normal variable:

$$\boldsymbol{\Sigma} = \mathbf{L}\mathbf{L}^T$$

- (3) compute  $\mathbf{x} = \mathbf{L}\mathbf{z} + \boldsymbol{\mu}$  – such  $\mathbf{x}$  is a sample from  $\mathcal{MVN}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$  distribution. ■

## (F)

*Sol:* We first note that a matrix quadratic form  $Q(\mathbf{x} - \boldsymbol{\mu})$  has the form

$$Q(\mathbf{x} - \boldsymbol{\mu}) = (\mathbf{x} - \boldsymbol{\mu})^T (\mathbf{A}(\mathbf{x} - \boldsymbol{\mu}))$$

in which  $\mathbf{A}$  is a symmetric matrix.

We've actually solved this problem in our solution for (D) under Method 1 – we applied change of variable on standard multivariate normal variable  $\mathbf{z}$  to derive the pdf of  $\mathbf{x} \sim \mathcal{MVN}(\boldsymbol{\mu}, \mathbf{L}\mathbf{L}^T)$ . We briefly go over the steps again here:

The pdf of the standard multivariate normal variable  $\mathbf{z}$  is

$$p(\mathbf{z}) = \frac{1}{(2\pi)^{p/2}} e^{-\mathbf{z}^T \mathbf{z}/2}$$

From (E) we have that  $\mathbf{x} \sim \mathcal{MVN}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$  can be expressed by  $\mathbf{z}$  as

$$\mathbf{x} = \mathbf{L}\mathbf{z} + \boldsymbol{\mu}$$

in which  $\mathbf{L}\mathbf{L}^T = \boldsymbol{\Sigma}$  is the Cholesky decompositon of  $\boldsymbol{\Sigma}$ .

Applying change of variable:

$$g(\mathbf{z}) = \mathbf{x} = \mathbf{L}\mathbf{z} + \boldsymbol{\mu}$$

note that  $\mathbf{L}$  here is invertible, thus we have

$$\mathbf{z} = g^{-1}(\mathbf{x}) = \mathbf{L}^{-1}(\mathbf{x} - \boldsymbol{\mu})$$

Therefore we have the Jacobian

$$J = \left| \frac{\partial \mathbf{r}(\mathbf{x})}{\partial \mathbf{x}} \right| = \left| \frac{\partial \mathbf{g}^{-1}(\mathbf{x})}{\partial \mathbf{x}} \right| = |\mathbf{L}^{-1}| = |\mathbf{L}|^{-1}$$

The pdf of  $\mathbf{x}$  is thus

$$\begin{aligned} p(\mathbf{x}) &= \frac{1}{(2\pi)^{p/2}} e^{-(\mathbf{L}^{-1}(\mathbf{x}-\boldsymbol{\mu}))^T(\mathbf{L}^{-1}(\mathbf{x}-\boldsymbol{\mu}))/2} \cdot |J| \\ &= (2\pi)^{-p/2} |\mathbf{L}|^{-1} \exp\left(-\frac{1}{2}(\mathbf{x}-\boldsymbol{\mu})^T(\mathbf{L}\mathbf{L}^T)^{-1}(\mathbf{x}-\boldsymbol{\mu})\right) \end{aligned}$$

Note that  $\mathbf{L}\mathbf{L}^T = \boldsymbol{\Sigma}$  which is a symmetric matrix, thus  $(\mathbf{x}-\boldsymbol{\mu})^T(\mathbf{L}\mathbf{L}^T)^{-1}(\mathbf{x}-\boldsymbol{\mu})$  is a quadratic form. We thus can express the pdf as

$$p(\mathbf{x}) = C \cdot \exp\left(-\frac{1}{2}Q(\mathbf{x}-\boldsymbol{\mu})\right)$$

in which

$$C = (2\pi)^{-p/2} |\mathbf{L}|^{-1}$$

and

$$Q(\mathbf{x}-\boldsymbol{\mu}) = (\mathbf{x}-\boldsymbol{\mu})^T(\boldsymbol{\Sigma})^{-1}(\mathbf{x}-\boldsymbol{\mu})$$

■

(G)

*Sol:* To solve this problem, we first compute the mgf of  $\mathbf{y} = \mathbf{Ax}_1 + \mathbf{Bx}_2$ , then use the fact that a mgf can uniquely determine a distribution to identify the distribution of  $\mathbf{y}$  with corresponding parameters.

We have the mgf of  $\mathbf{y}$  as

$$\begin{aligned}
M_{\mathbf{y}}(\mathbf{t}) &= \mathbb{E}(e^{\mathbf{t}^T \mathbf{y}}) \\
&= \mathbb{E}(e^{\mathbf{t}^T (\mathbf{A}\mathbf{x}_1 + \mathbf{B}\mathbf{x}_2)}) \\
&= \mathbb{E}(e^{\mathbf{t}^T (\mathbf{A}\mathbf{x}_1)} e^{\mathbf{t}^T (\mathbf{B}\mathbf{x}_2)}) \\
&= \mathbb{E}(e^{(\mathbf{A}^T \mathbf{t})^T \mathbf{x}_1} e^{(\mathbf{B}^T \mathbf{t})^T \mathbf{x}_2}) \\
&= \int_{\mathbf{x}_1} \int_{\mathbf{x}_2} e^{(\mathbf{A}^T \mathbf{t})^T \mathbf{x}_1} e^{(\mathbf{B}^T \mathbf{t})^T \mathbf{x}_2} f(\mathbf{x}_1, \mathbf{x}_2) d\mathbf{x}_1 d\mathbf{x}_2 \\
&= \int_{\mathbf{x}_1} \int_{\mathbf{x}_2} e^{(\mathbf{A}^T \mathbf{t})^T \mathbf{x}_1} e^{(\mathbf{B}^T \mathbf{t})^T \mathbf{x}_2} f(\mathbf{x}_1) f(\mathbf{x}_2) d\mathbf{x}_1 d\mathbf{x}_2 \quad (\mathbf{x}_1 \text{ and } \mathbf{x}_2 \text{ are independent}) \\
&= \int_{\mathbf{x}_1} e^{(\mathbf{A}^T \mathbf{t})^T \mathbf{x}_1} f(\mathbf{x}_1) d\mathbf{x}_1 \cdot \int_{\mathbf{x}_2} e^{(\mathbf{B}^T \mathbf{t})^T \mathbf{x}_2} f(\mathbf{x}_2) d\mathbf{x}_2 \\
&= M_{\mathbf{x}_1}(\mathbf{A}^T \mathbf{t}) \cdot M_{\mathbf{x}_2}(\mathbf{B}^T \mathbf{t}) \\
&= \exp\left((\mathbf{A}^T \mathbf{t})^T \boldsymbol{\mu}_1 + \frac{(\mathbf{A}^T \mathbf{t})^T \boldsymbol{\Sigma}_1 (\mathbf{A}^T \mathbf{t})}{2}\right) \cdot \exp\left((\mathbf{B}^T \mathbf{t})^T \boldsymbol{\mu}_2 + \frac{(\mathbf{B}^T \mathbf{t})^T \boldsymbol{\Sigma}_2 (\mathbf{B}^T \mathbf{t})}{2}\right) \\
&= \exp\left(\mathbf{t}^T (\mathbf{A}\boldsymbol{\mu}_1 + \mathbf{B}\boldsymbol{\mu}_2) + \frac{\mathbf{t}^T (\mathbf{A}\boldsymbol{\Sigma}_1 \mathbf{A}^T + \mathbf{B}\boldsymbol{\Sigma}_2 \mathbf{B}^T) \mathbf{t}}{2}\right)
\end{aligned}$$

This is almost the mgf of a multivariate normal variable, except that we need to justify that  $(\mathbf{A}\boldsymbol{\Sigma}_1 \mathbf{A}^T + \mathbf{B}\boldsymbol{\Sigma}_2 \mathbf{B}^T)$  is positive definite so that it can be a covariance matrix:

Since both  $\mathbf{A}$  and  $\mathbf{B}$  have full column rank, and both  $\boldsymbol{\Sigma}_1$  and  $\boldsymbol{\Sigma}_2$  are positive definite since they are both covariance matrices,  $\mathbf{A}\boldsymbol{\Sigma}_1 \mathbf{A}^T$  and  $\mathbf{B}\boldsymbol{\Sigma}_2 \mathbf{B}^T$  are each positive definite. The sum of two positive definite matrices (with same dimension) is positive definite too: denote two p.d. matrices as  $\mathbf{M}$  and  $\mathbf{N}$ , s.t.

$$\mathbf{u}^T \mathbf{M} \mathbf{u} > 0$$

and

$$\mathbf{u}^T \mathbf{N} \mathbf{u} > 0$$

for all nonzero vectors  $\mathbf{u}$ , then

$$\mathbf{u}^T (\mathbf{M} + \mathbf{N}) \mathbf{u} = \mathbf{u}^T \mathbf{M} \mathbf{u} + \mathbf{u}^T \mathbf{N} \mathbf{u} > 0$$

Therefore,  $(\mathbf{A}\boldsymbol{\Sigma}_1 \mathbf{A}^T + \mathbf{B}\boldsymbol{\Sigma}_2 \mathbf{B}^T)$  is positive definite, thus eligible to be a covariance matrix. We thus have

$$\mathbf{y} \sim \mathcal{MVN}\left(\mathbf{A}\boldsymbol{\mu}_1 + \mathbf{B}\boldsymbol{\mu}_2, \mathbf{A}\boldsymbol{\Sigma}_1 \mathbf{A}^T + \mathbf{B}\boldsymbol{\Sigma}_2 \mathbf{B}^T\right)$$

■

## Conditionals and Marginals

(A)

*Sol:*

## Method 1

This solution applies the affine transformation of  $\mathbf{x}$ : previously, we write  $\mathbf{x}$  as an affine transformation of a standard multivariate normal random variable; this time, we directly apply affine transformation on  $\mathbf{x}$ .

Define  $\mathbf{y} = \mathbf{Ax}$  in which

$$\mathbf{A} = [\mathbf{I}_{k \times k} \quad \mathbf{0}_{(p-k) \times k}]$$

thus a  $k$ -by- $p$  matrix with full row rank.

For  $\mathbf{x} \sim \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ , we have the mgf of  $\mathbf{y}$  as

$$M_{\mathbf{y}}(\mathbf{t}) = \mathbb{E}(e^{\mathbf{t}^T \mathbf{Ax}}) = \mathbb{E}(e^{(\mathbf{A}^T \mathbf{t})^T \mathbf{x}}) = \exp((\mathbf{A}^T \mathbf{t})^T \boldsymbol{\mu}, \frac{(\mathbf{A}^T \mathbf{t})^T \boldsymbol{\Sigma} (\mathbf{A}^T \mathbf{t})}{2}) = \exp(\mathbf{t}^T (\mathbf{A} \boldsymbol{\mu}) + \frac{\mathbf{t}^T (\mathbf{A} \boldsymbol{\Sigma} \mathbf{A}^T) \mathbf{t}}{2})$$

Since  $\mathbf{A}^T$  has full column rank, given any non-zero vector  $\mathbf{v}$  we have  $\mathbf{A}^T \mathbf{v}$  as a non-zero vector. Since  $\boldsymbol{\Sigma}$  as a covariance matrix is positive definite, we have

$$\mathbf{v}^T (\mathbf{A} \boldsymbol{\Sigma} \mathbf{A}^T) \mathbf{v} = (\mathbf{A}^T \mathbf{v})^T \boldsymbol{\Sigma} (\mathbf{A}^T \mathbf{v}) > 0$$

therefore  $\mathbf{A} \boldsymbol{\Sigma} \mathbf{A}^T$  is positive definite, thus eligible to be a covariance matrix. By our conclusion from Part (C) of previous question,

$$\mathbf{y} \sim \mathcal{MVN}(\mathbf{A}\boldsymbol{\mu}, \mathbf{A}\boldsymbol{\Sigma}\mathbf{A}^T)$$

Using the notation in our question, we have  $\mathbf{y} = \mathbf{Ax} = \mathbf{x}_1$ ,  $\mathbf{A}\boldsymbol{\mu} = \boldsymbol{\mu}_1$ ,  $\mathbf{A}\boldsymbol{\Sigma}\mathbf{A}^T = \boldsymbol{\Sigma}_{11}$ , therefore in other words,

$$\mathbf{x}_1 \sim \mathcal{MVN}(\boldsymbol{\mu}_1, \boldsymbol{\Sigma}_{11})$$

## Method 2

This method uses the affine transformation of standard multivariate normal variable that we studied from the previous question. I learnt this method from class and like it very much:

For  $\mathbf{x} = \mathbf{Lz} + \boldsymbol{\mu}$  in which

$$\mathbf{L} = \begin{bmatrix} \mathbf{L}_{11} & \mathbf{L}_{12} \\ \mathbf{L}_{21} & \mathbf{L}_{22} \end{bmatrix}$$

we have

$$\begin{aligned} \mathbf{L} \cdot \mathbf{L}^T &= \begin{bmatrix} \mathbf{L}_{11} & \mathbf{L}_{12} \\ \mathbf{L}_{21} & \mathbf{L}_{22} \end{bmatrix} \cdot \begin{bmatrix} \mathbf{L}_{11}^T & \mathbf{L}_{21}^T \\ \mathbf{L}_{12}^T & \mathbf{L}_{22}^T \end{bmatrix} \\ &= \begin{bmatrix} \mathbf{L}_{11} \cdot \mathbf{L}_{11}^T + \mathbf{L}_{12} \cdot \mathbf{L}_{12}^T & \mathbf{L}_{11} \cdot \mathbf{L}_{21}^T + \mathbf{L}_{12} \cdot \mathbf{L}_{22}^T \\ \mathbf{L}_{21} \cdot \mathbf{L}_{11}^T + \mathbf{L}_{22} \cdot \mathbf{L}_{12}^T & \mathbf{L}_{21} \cdot \mathbf{L}_{21}^T + \mathbf{L}_{22} \cdot \mathbf{L}_{22}^T \end{bmatrix} \\ &= \boldsymbol{\Sigma} \text{ (we have } \mathbf{L} \mathbf{L}^T = \boldsymbol{\Sigma} \text{ from Part (E) of previous question)} \\ &= \begin{bmatrix} \boldsymbol{\Sigma}_{11} & \boldsymbol{\Sigma}_{12} \\ \boldsymbol{\Sigma}_{21} & \boldsymbol{\Sigma}_{22} \end{bmatrix} \end{aligned}$$

Meanwhile,

$$\begin{aligned}
\begin{bmatrix} \mathbf{x}_1 \\ \mathbf{x}_2 \end{bmatrix} &= \mathbf{x} \\
&= \mathbf{L}\mathbf{z} + \boldsymbol{\mu} \\
&= \begin{bmatrix} \mathbf{L}_{11} & \mathbf{L}_{12} \\ \mathbf{L}_{21} & \mathbf{L}_{22} \end{bmatrix} \begin{bmatrix} \mathbf{z}_1 \\ \mathbf{z}_2 \end{bmatrix} + \begin{bmatrix} \boldsymbol{\mu}_1 \\ \boldsymbol{\mu}_2 \end{bmatrix} \\
&= \begin{bmatrix} \mathbf{L}_{11}\mathbf{z}_1 + \mathbf{L}_{12}\mathbf{z}_2 \\ \mathbf{L}_{21}\mathbf{z}_1 + \mathbf{L}_{22}\mathbf{z}_2 \end{bmatrix} + \begin{bmatrix} \boldsymbol{\mu}_1 \\ \boldsymbol{\mu}_2 \end{bmatrix}
\end{aligned}$$

in which

$$\mathbf{x}_1 = \mathbf{L}_{11}\mathbf{z}_1 + \mathbf{L}_{12}\mathbf{z}_2 + \boldsymbol{\mu}_1$$

therefore, the mgf of  $\mathbf{x}_1$  is

$$\begin{aligned}
M_{\mathbf{x}_1}(\mathbf{t}) &= \mathbb{E}[\exp(\mathbf{t}^T(\mathbf{L}_{11}\mathbf{z}_1 + \mathbf{L}_{12}\mathbf{z}_2 + \boldsymbol{\mu}_1))] \\
&= \mathbb{E}[\exp(\mathbf{L}_{11}^T\mathbf{t})^T\mathbf{z}_1] \cdot \mathbb{E}[\exp(\mathbf{L}_{12}^T\mathbf{t})^T\mathbf{z}_2] \cdot \exp(\mathbf{t}^T\boldsymbol{\mu}_1) \\
&= \exp\left(\frac{1}{2}(\mathbf{L}_{11}^T\mathbf{t})^T(\mathbf{L}_{11}^T\mathbf{t})\right) \exp\left(\frac{1}{2}(\mathbf{L}_{12}^T\mathbf{t})^T(\mathbf{L}_{12}^T\mathbf{t})\right) \exp(\mathbf{t}^T\boldsymbol{\mu}_1) \\
&\quad (\mathbf{z}_1 \text{ and } \mathbf{z}_2 \text{ are independent standard multivariate normal variables}) \\
&= \exp\left(\mathbf{t}^T\boldsymbol{\mu}_1 + \frac{1}{2}\mathbf{t}^T(\mathbf{L}_{11} \cdot \mathbf{L}_{11}^T + \mathbf{L}_{12} \cdot \mathbf{L}_{12}^T)\mathbf{t}\right) \\
&= \exp(\mathbf{t}^T\boldsymbol{\mu}_1 + \frac{1}{2}\mathbf{t}^T\Sigma_{11}\mathbf{t})
\end{aligned}$$

Here we are very tempted to write out the distribution of  $\mathbf{x}_1$ , but we do need to argue that  $\Sigma_{11}$  is positive definite. My current solution is to construct  $k$ -by- $p$  matrix  $\mathbf{A} = [\mathbf{I}_{k \times k} \ \mathbf{0}_{(p-k) \times k}]$  as in **Method 1**, s.t.  $\Sigma_{11} = \mathbf{A}\boldsymbol{\Sigma}\mathbf{A}^T$  being positive definite is what we've showed before; but I haven't figured out a different way that circumvents such construction of  $\mathbf{A}$ .

After the above arguments, we can now write that

$$\mathbf{x}_1 \sim \mathcal{MVN}(\boldsymbol{\mu}_1, \Sigma_{11})$$

by the form of mgf. ■

(B)

*Sol:* Before we conduct the derivation, I just want to point out that  $\Omega$ , just like  $\Sigma$ , is symmetric:

Given any symmetric matrix  $\mathbf{A}$ , we have

$$\begin{aligned}
\mathbf{A}^{-1}\mathbf{A} &= \mathbf{I}_{p \times p} \\
\Rightarrow (\mathbf{A}^{-1}\mathbf{A})^T &= \mathbf{A}^T(\mathbf{A}^{-1})^T = \mathbf{I} \\
\Rightarrow \mathbf{A}(\mathbf{A}^{-1})^T &= \mathbf{I} \\
\Rightarrow (\mathbf{A}^{-1})^T &= \mathbf{A}^{-1}
\end{aligned}$$

Therefore, we might write  $\Omega$  as

$$\Omega = \begin{bmatrix} \Omega_{11} & \Omega_{12} \\ \Omega_{12}^T & \Omega_{22} \end{bmatrix}$$

Since  $\Omega = \Sigma^{-1}$ , we have

$$\begin{aligned}
& \Omega \Sigma = I_{p \times p} \\
& \Rightarrow \begin{bmatrix} \Omega_{11} & \Omega_{12} \\ \Omega_{12}^T & \Omega_{22} \end{bmatrix} \cdot \begin{bmatrix} \Sigma_{11} & \Sigma_{12} \\ \Sigma_{12}^T & \Sigma_{22} \end{bmatrix} = \begin{bmatrix} I_{k \times k} & \mathbf{0}_{k \times (p-k)} \\ \mathbf{0}_{(p-k) \times k} & I_{(p-k) \times (p-k)} \end{bmatrix} \\
& \Rightarrow \begin{bmatrix} \Omega_{11}\Sigma_{11} + \Omega_{12}\Sigma_{12}^T & \Omega_{11}\Sigma_{12} + \Omega_{12}\Sigma_{22} \\ \Omega_{12}^T\Sigma_{11} + \Omega_{22}\Sigma_{12}^T & \Omega_{12}^T\Sigma_{12} + \Omega_{22}\Sigma_{22} \end{bmatrix} = \begin{bmatrix} I_{k \times k} & \mathbf{0}_{k \times (p-k)} \\ \mathbf{0}_{(p-k) \times k} & I_{(p-k) \times (p-k)} \end{bmatrix} \\
& \Rightarrow \begin{cases} \Omega_{11}\Sigma_{11} + \Omega_{12}\Sigma_{12}^T = I_{k \times k} & (1) \\ \Omega_{11}\Sigma_{12} + \Omega_{12}\Sigma_{22} = \mathbf{0}_{k \times (p-k)} & (2) \\ \Omega_{12}^T\Sigma_{11} + \Omega_{22}\Sigma_{12}^T = \mathbf{0}_{(p-k) \times k} & (3) \\ \Omega_{12}^T\Sigma_{12} + \Omega_{22}\Sigma_{22} = I_{(p-k) \times (p-k)} & (4) \end{cases}
\end{aligned}$$

From Equation (2) we have

$$\begin{aligned}
& \Omega_{12}\Sigma_{22} = -\Omega_{11}\Sigma_{12} \\
& \Rightarrow \Omega_{12} = -\Omega_{11}\Sigma_{12}\Sigma_{22}^{-1}
\end{aligned}$$

Plugging it into Equation (1) we have

$$\begin{aligned}
& \Omega_{11}\Sigma_{11} - \Omega_{11}\Sigma_{12}\Sigma_{22}^{-1}\Sigma_{12}^T = I_{k \times k} \\
& \Rightarrow \Omega_{11}(\Sigma_{11} - \Sigma_{12}\Sigma_{22}^{-1}\Sigma_{12}^T) = I_{k \times k} \\
& \Rightarrow \Omega_{11} = (\Sigma_{11} - \Sigma_{12}\Sigma_{22}^{-1}\Sigma_{12}^T)^{-1}
\end{aligned}$$

thus

$$\Omega_{12} = -(\Sigma_{11} - \Sigma_{12}\Sigma_{22}^{-1}\Sigma_{12}^T)^{-1}\Sigma_{12}\Sigma_{22}^{-1}$$

Meanwhile, with similar arguments, we have from Equation (3) s.t.

$$\Omega_{12}^T = -\Omega_{22}\Sigma_{12}^T\Sigma_{11}^{-1}$$

Plugging it into Equation (4) we have

$$\begin{aligned}
& -\Omega_{22}\Sigma_{12}^T\Sigma_{11}^{-1}\Sigma_{12} + \Omega_{22}\Sigma_{22} = I_{(p-k) \times (p-k)} \\
& \Rightarrow \Omega_{22}(\Sigma_{22} - \Sigma_{12}^T\Sigma_{11}^{-1}\Sigma_{12}) = I_{(p-k) \times (p-k)} \\
& \Rightarrow \Omega_{22} = (\Sigma_{22} - \Sigma_{12}^T\Sigma_{11}^{-1}\Sigma_{12})^{-1}
\end{aligned}$$

and another expression for  $\omega_{12}^T$  (although we've derived  $\omega_{12}$  above):

$$\Omega_{12}^T = -(\Sigma_{22} - \Sigma_{12}^T\Sigma_{11}^{-1}\Sigma_{12})^{-1}\Sigma_{12}^T\Sigma_{11}^{-1}$$

We thus re-write  $\Omega$  as follows, in terms of blocks of  $\Sigma$ :

$$\Omega = \begin{bmatrix} \Omega_{11} & \Omega_{12} \\ \Omega_{12}^T & \Omega_{22} \end{bmatrix} = \begin{bmatrix} (\Sigma_{11} - \Sigma_{12}\Sigma_{22}^{-1}\Sigma_{12}^T)^{-1} & -(\Sigma_{11} - \Sigma_{12}\Sigma_{22}^{-1}\Sigma_{12}^T)^{-1}\Sigma_{12}\Sigma_{22}^{-1} \\ -(\Sigma_{22} - \Sigma_{12}^T\Sigma_{11}^{-1}\Sigma_{12})^{-1}\Sigma_{12}^T\Sigma_{11}^{-1} & (\Sigma_{22} - \Sigma_{12}^T\Sigma_{11}^{-1}\Sigma_{12})^{-1} \end{bmatrix}$$

■

## (C)

*Sol:* For this question, we first derive the conditional distribution in terms of the blocks of precision matrix (along with  $\mathbf{x}$  and  $\boldsymbol{\mu}$ ), then convert the parameters into terms of covariance matrix.

Now we begin the derivation:

to derive the conditional distribution of  $\mathbf{x}_1$ , we could utilize the relationship indicated by the following equation to derive the analytic form of the conditional pdf:

$$p(\mathbf{x}) = p(\mathbf{x}_1, \mathbf{x}_2) = p(\mathbf{x}_1 | \mathbf{x}_2) \cdot p(\mathbf{x}_2)$$

in which we could go through similar arguments as in (A) to derive the parameters of the marginal distribution of  $x_2$ , we don't need to do that – we can simply treat  $\mathbf{x}_2$  as constant, and keep simplifying the pdf form until we only have terms including  $\mathbf{x}_1$ , while other terms would just be absorbed into the normalizing constant; and clearly, the marginal pdf of  $\mathbf{x}_2$ ,  $p(\mathbf{x}_2)$ , won't have the term  $\mathbf{x}_1$ .

Since the term  $\mathbf{x}_2$  only appears in the exponent of the base of natural logarithm, we take the question's advice to work with the densities on a log scale, and gradually removing the terms that can be treated as constants. Again, notice that we are working with precision matrix during the derivation, instead of covariance matrix:

$$\begin{aligned} \log p(\mathbf{x}_1 | \mathbf{x}_2) &\propto \log p(\mathbf{x}_1, \mathbf{x}_2) \\ &\propto \log \exp \left( -\frac{1}{2} \left( \begin{pmatrix} \mathbf{x}_1 \\ \mathbf{x}_2 \end{pmatrix} - \begin{pmatrix} \boldsymbol{\mu}_1 \\ \boldsymbol{\mu}_2 \end{pmatrix} \right)^T \begin{pmatrix} \boldsymbol{\Omega}_{11} & \boldsymbol{\Omega}_{12} \\ \boldsymbol{\Omega}_{12}^T & \boldsymbol{\Omega}_{22} \end{pmatrix} \left( \begin{pmatrix} \mathbf{x}_1 \\ \mathbf{x}_2 \end{pmatrix} - \begin{pmatrix} \boldsymbol{\mu}_1 \\ \boldsymbol{\mu}_2 \end{pmatrix} \right) \right) \\ &\propto \left( \begin{pmatrix} \mathbf{x}_1 \\ \mathbf{x}_2 \end{pmatrix} - \begin{pmatrix} \boldsymbol{\mu}_1 \\ \boldsymbol{\mu}_2 \end{pmatrix} \right)^T \begin{pmatrix} \boldsymbol{\Omega}_{11} & \boldsymbol{\Omega}_{12} \\ \boldsymbol{\Omega}_{12}^T & \boldsymbol{\Omega}_{22} \end{pmatrix} \left( \begin{pmatrix} \mathbf{x}_1 \\ \mathbf{x}_2 \end{pmatrix} - \begin{pmatrix} \boldsymbol{\mu}_1 \\ \boldsymbol{\mu}_2 \end{pmatrix} \right) \\ &= (\mathbf{x}_1 - \boldsymbol{\mu}_1)^T \boldsymbol{\Omega}_{11} (\mathbf{x}_1 - \boldsymbol{\mu}_1) + (\mathbf{x}_2 - \boldsymbol{\mu}_2)^T \boldsymbol{\Omega}_{12}^T (\mathbf{x}_1 - \boldsymbol{\mu}_1) \\ &\quad + (\mathbf{x}_1 - \boldsymbol{\mu}_1)^T \boldsymbol{\Omega}_{12} (\mathbf{x}_2 - \boldsymbol{\mu}_2) + (\mathbf{x}_2 - \boldsymbol{\mu}_2)^T \boldsymbol{\Omega}_{22} (\mathbf{x}_2 - \boldsymbol{\mu}_2) \\ &\propto (\mathbf{x}_1 - \boldsymbol{\mu}_1)^T \boldsymbol{\Omega}_{11} (\mathbf{x}_1 - \boldsymbol{\mu}_1) + (\mathbf{x}_2 - \boldsymbol{\mu}_2)^T \boldsymbol{\Omega}_{12}^T (\mathbf{x}_1 - \boldsymbol{\mu}_1) \\ &\quad + (\mathbf{x}_1 - \boldsymbol{\mu}_1)^T \boldsymbol{\Omega}_{12} (\mathbf{x}_2 - \boldsymbol{\mu}_2) \\ &= (\mathbf{x}_1 - \boldsymbol{\mu}_1)^T \boldsymbol{\Omega}_{11} (\mathbf{x}_1 - \boldsymbol{\mu}_1) + (\mathbf{x}_2 - \boldsymbol{\mu}_2)^T \boldsymbol{\Omega}_{12}^T \boldsymbol{\Omega}_{11}^{-1} (\mathbf{x}_1 - \boldsymbol{\mu}_1) \\ &\quad + (\mathbf{x}_1 - \boldsymbol{\mu}_1)^T \boldsymbol{\Omega}_{11} \boldsymbol{\Omega}_{11}^{-1} \boldsymbol{\Omega}_{12} (\mathbf{x}_2 - \boldsymbol{\mu}_2) \\ &\propto ((\mathbf{x}_1 - \boldsymbol{\mu}_1) + \boldsymbol{\Omega}_{11}^{-1} \boldsymbol{\Omega}_{12} (\mathbf{x}_2 - \boldsymbol{\mu}_2))^T \boldsymbol{\Omega}_{11} ((\mathbf{x}_1 - \boldsymbol{\mu}_1) + \boldsymbol{\Omega}_{11}^{-1} \boldsymbol{\Omega}_{12} (\mathbf{x}_2 - \boldsymbol{\mu}_2)) \end{aligned}$$

Converting above result back from log scale, we have

$$\begin{aligned} p(\mathbf{x}_1 | \mathbf{x}_2) &= \exp(\log p(\mathbf{x}_1 | \mathbf{x}_2)) \\ &\propto \exp \left( -\frac{1}{2} ((\mathbf{x}_1 - \boldsymbol{\mu}_1) + \boldsymbol{\Omega}_{11}^{-1} \boldsymbol{\Omega}_{12} (\mathbf{x}_2 - \boldsymbol{\mu}_2))^T \boldsymbol{\Omega}_{11} ((\mathbf{x}_1 - \boldsymbol{\mu}_1) + \boldsymbol{\Omega}_{11}^{-1} \boldsymbol{\Omega}_{12} (\mathbf{x}_2 - \boldsymbol{\mu}_2)) \right) \end{aligned}$$

which is exactly the form of the kernel of  $\mathcal{MVN}(\boldsymbol{\mu}_1 - \boldsymbol{\Omega}_{11}^{-1} \boldsymbol{\Omega}_{12} (\mathbf{x}_2 - \boldsymbol{\mu}_2), \boldsymbol{\Omega}_{11}^{-1})$

We temporarily leave out the discussion on the positive-definiteness of  $\boldsymbol{\Omega}_{11}^{-1}$ , as it requires more linear-algebraic arguments; for now, I simply point out that it's symmetric:

$$(\boldsymbol{\Omega}_{11}^{-1})^T = (\boldsymbol{\Sigma}_{11} - \boldsymbol{\Sigma}_{12} \boldsymbol{\Sigma}_{22}^{-1} \boldsymbol{\Sigma}_{12}^T)^T = \boldsymbol{\Sigma}_{11}^T - \boldsymbol{\Sigma}_{12} (\boldsymbol{\Sigma}_{22}^{-1})^T \boldsymbol{\Sigma}_{12}^T = \boldsymbol{\Sigma}_{11}^T - \boldsymbol{\Sigma}_{12} (\boldsymbol{\Sigma}_{22}^T)^{-1} \boldsymbol{\Sigma}_{12}^T = \boldsymbol{\Sigma}_{11} - \boldsymbol{\Sigma}_{12} \boldsymbol{\Sigma}_{22}^{-1} \boldsymbol{\Sigma}_{12}^T = \boldsymbol{\Omega}_{11}^{-1}$$

Therefore, in terms of (blocks of) covariance matrix, by our derivation in (B) we have we thus have the conditional distribution

$$\mathbf{x}_1 | \mathbf{x}_2 \sim \mathcal{MVN}(\boldsymbol{\mu}_1 + \boldsymbol{\Sigma}_{12} \boldsymbol{\Sigma}_{22}^{-1} (\mathbf{x}_2 - \boldsymbol{\mu}_2), \boldsymbol{\Sigma}_{11} - \boldsymbol{\Sigma}_{12} \boldsymbol{\Sigma}_{22}^{-1} \boldsymbol{\Sigma}_{12}^T)$$

■

# Multiple Regression: Three Classical Principles for Inference

(A)

*Sol:* We show the derivation from each of the three principles individually as follows:

## Least Squares

To find the minimum of the sum of squares w.r.t.  $\beta$ , we first compute the critical point of the sum by solving for the  $\beta$  when its gradient equals to 0, then justifying the sum reaches its minimum by showing the Hessian matrix of  $\beta$  is positive definite:

For

$$\begin{aligned}\sum(y_i - \mathbf{x}_i^T \beta)^2 &= \sum(y_i - \mathbf{x}_i^T \beta)^T (y_i - \mathbf{x}_i^T \beta) \\ &= \sum(y_i^T - \beta^T \mathbf{x}_i)(y_i - \mathbf{x}_i^T \beta) \\ &= \sum(y_i^T y_i - 2y_i^T \mathbf{x}_i^T \beta + \beta^T \mathbf{x}_i \mathbf{x}_i^T \beta)\end{aligned}$$

Note that although each  $y_i$  is a scalar, we treat them as a 1-by-1 matrix, thus has a meaning of being transposed. The gradient w.r.t.  $\beta$  is thus

$$\frac{\partial}{\partial \beta} \sum(y_i^T y_i - 2y_i^T \mathbf{x}_i^T \beta + \beta^T \mathbf{x}_i \mathbf{x}_i^T \beta) = -2(\sum \mathbf{x}_i y_i) + 2(\sum \mathbf{x}_i \mathbf{x}_i^T \beta)$$

Setting it equal to  $\mathbf{0}$ , we have

$$\begin{aligned}2(\sum \mathbf{x}_i y_i) + 2(\sum \mathbf{x}_i \mathbf{x}_i^T \beta^*) &= \mathbf{0} \\ \Rightarrow \sum \mathbf{x}_i \mathbf{x}_i^T \beta^* &= \sum \mathbf{x}_i y_i\end{aligned}\tag{1}$$

It's tempting to leave it like that and derive for  $\beta^*$ , but we rearrange the terms to make the summation sign disappear, since it would be much easier this way to argue that the Hessian matrix (which we'd compute later) is positive definite.

We define

$$\mathbf{X} = \begin{bmatrix} \mathbf{x}_1^T \\ \mathbf{x}_2^T \\ \vdots \\ \mathbf{x}_n^T \end{bmatrix} \Rightarrow \mathbf{X}^T = [\mathbf{x}_1 \quad \mathbf{x}_2 \quad \dots \quad \mathbf{x}_n], \mathbf{Y} = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix}$$

we can see

$$\sum \mathbf{x}_i \mathbf{x}_i^T = \mathbf{X}^T \mathbf{X}, \sum \mathbf{x}_i y_i = \mathbf{X}^T \mathbf{Y}$$

We shall directly claim that  $\mathbf{X}$  has full column rank, because each of its column can be viewed as a feature, and generally all of the features shall be linearly independent – otherwise why bother having the extra feature that's linearly dependent on others, since we want the feature to help provide extra information in predicting  $\mathbf{Y}$ . Therefore, for  $n$ -by- $p$  matrix  $\mathbf{X}$ , given any  $p$ -by-1 vector  $\mathbf{u}$ , we shall have  $\mathbf{X}\mathbf{u}$  as a non-zero  $n$ -by-1 vector; or in other words, the null space of  $\mathbf{X}$  is  $\{\mathbf{0}\}$ . We thus have

$$\mathbf{u}^T \mathbf{X}^T \mathbf{X} \mathbf{u} = (\mathbf{X}\mathbf{u})^T (\mathbf{X}\mathbf{u}) > 0$$

i.e.  $\mathbf{X}^T \mathbf{X}$  is positive definite. Therefore, it is invertible, we can thus have  $\beta^* = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{Y}$

To argue that such  $\beta^*$  gives us the minimum of  $\sum(y_i - \mathbf{x}_i^T \beta)^2$ , we compute the Hessian matrix of  $\beta$ :

$$\mathcal{H}(\boldsymbol{\beta}) = \frac{\partial}{\partial \boldsymbol{\beta}} \frac{\partial \sum (y_i - \mathbf{x}_i^T \boldsymbol{\beta})^2}{\partial \boldsymbol{\beta}} = 2 \sum \mathbf{x}_i \mathbf{x}_i^T = 2 \mathbf{X}^T \mathbf{X}$$

which is clearly positive definite by our conclusion from above. Therefore, we'd achieve the minimum of  $\sum (y_i - \mathbf{x}_i^T \boldsymbol{\beta})^2$  at

$$\hat{\boldsymbol{\beta}} = \boldsymbol{\beta}^* = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{Y}$$

### Maximum Likelihood under Gaussianity

For this question, we want to first find out the distribution for each  $y_i$ . After we do that, we plug the corresponding pdf into the function to be optimized. Eventually, we'd show that this optimization problem is actually equivalent to the one from **Least Squares**, thus can get the result for free by referring to the derivation from last section.

We first find out each  $y_i$ 's distribution by computing its mgf: for each  $y_i = \mathbf{x}_i^T \boldsymbol{\beta} + \epsilon_i$ , in which  $\epsilon_i \sim \mathcal{N}(0, \sigma^2)$ , we have

$$\begin{aligned} M_{y_i}(t) &= \mathbb{E}(e^{y_i t}) \\ &= \mathbb{E}(e^{(\mathbf{x}_i^T \boldsymbol{\beta} + \epsilon_i)t}) \\ &= e^{(\mathbf{x}_i^T \boldsymbol{\beta})t} \cdot \mathbb{E}(e^{\epsilon_i t}) \quad (\mathbf{x}_i^T \boldsymbol{\beta} \text{ is a constant term}) \\ &= e^{(\mathbf{x}_i^T \boldsymbol{\beta})t} \cdot e^{\sigma^2 t^2 / 2} \\ &= \exp((\mathbf{x}_i^T \boldsymbol{\beta})t + \sigma^2 t^2 / 2) \end{aligned}$$

We recognize that the last expression is the mgf of a  $\mathcal{N}(\mathbf{x}_i^T \boldsymbol{\beta}, \sigma^2)$  random variable, thus is the distribution of each  $y_i$ .

Now using the fact that  $y_i$ 's are independent, we have their joint distribution

$$p(y_1, \dots, y_n) = \prod_{i=1}^n p(y_i) = \prod_{i=1}^n \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(y_i - \mathbf{x}_i^T \boldsymbol{\beta})^2}{2\sigma^2}\right)$$

To maximize this joint likelihood, we have the following equivalent statements:

$$\begin{aligned} &\arg \max_{\boldsymbol{\beta}} \left\{ \prod_{i=1}^n p(y_i) \right\} \\ &\equiv \arg \max_{\boldsymbol{\beta}} \left\{ \log \prod_{i=1}^n p(y_i) \right\} \\ &\equiv \arg \max_{\boldsymbol{\beta}} \left\{ \sum_{i=1}^n \log p(y_i) \right\} \\ &\equiv \arg \max_{\boldsymbol{\beta}} \left\{ \sum_{i=1}^n \left( -\frac{(y_i - \mathbf{x}_i^T \boldsymbol{\beta})^2}{2\sigma^2} \right) \right\} \quad (\text{removing constants}) \\ &\equiv \arg \min_{\boldsymbol{\beta}} \left\{ \sum_{i=1}^n ((y_i - \mathbf{x}_i^T \boldsymbol{\beta})^2) \right\} \quad (\text{flip negative sign to positive}) \end{aligned}$$

The very last expression shows that, to find  $\hat{\boldsymbol{\beta}}$  s.t. it maximizes the joint likelihood of  $y_i$ 's is equivalent to find the  $\hat{\boldsymbol{\beta}}$  that minimize the sum of squared errors. Therefore, the rest of derivation would just follow the same pattern as previous section, and of course we'd reach the same result.

## Method of Moments

For this section, we first show that centering the data won't change the covariance: for our random variable  $\mathbf{x}_i$ , denote centered variable as  $\mathbf{z}_i = \mathbf{x}_i - \mathbb{E}(\mathbf{x}_i)$ , such that  $\mathbb{E}(z_i) = \mathbb{E}(\mathbf{x}_i - \mathbb{E}(\mathbf{x}_i)) = \mathbf{0}$ . For any variable  $\mathbf{a}$ , we have

$$\text{cov}(\mathbf{z}_i, \mathbf{a}) = \mathbb{E}[(\mathbf{z}_i - \mathbb{E}(\mathbf{z}_i))(\mathbf{a} - \mathbb{E}(\mathbf{a}))] = \mathbb{E}[(\mathbf{z}_i)(\mathbf{a} - \mathbb{E}(\mathbf{a}))] = \mathbb{E}[(\mathbf{x}_i - \mathbb{E}(\mathbf{x}_i))(\mathbf{a} - \mathbb{E}(\mathbf{a}))] = \text{cov}(\mathbf{x}_i, \mathbf{a})$$

Therefore, we directly work with centered data. Assuming that our data is already centered, so that we can continue to use our previous notation  $\mathbf{x}_i$ . For each of the  $p$  predictors in random variable  $\mathbf{x}_i$ , we have the covariance between it and error term as

$$\text{cov}(x_{ij}, \epsilon_i) = \mathbb{E}[(x_{ij} - \mathbb{E}[x_{ij}])(\epsilon_i - \mathbb{E}[\epsilon_i])] = \mathbb{E}[(x_{ij})(\epsilon_i - \mathbb{E}[\epsilon_i])] = \mathbb{E}[x_{ij} \cdot \epsilon_i]$$

in which the last equation is true as we assume  $\mathbb{E}[\epsilon_i] = 0$

Applying method of moments, we use sample moments to estimate the population moments. We have the sample covariance between  $j$ -th predictor and error term  $Q_j$  as

$$Q_j = \frac{1}{n-1} \sum_{i=1}^n (x_{ij} - \bar{x}_j) \cdot \epsilon_i$$

in which the sample means are  $\bar{x}_j = \frac{1}{n} \sum_{i=1}^n x_{ij} = 0$ . This equation is true because the data is centered.

By setting the sample covariance to 0, we have

$$\sum_{i=1}^n x_{ij} \cdot \epsilon_i = 0 \Rightarrow \sum_{i=1}^n x_{ij} \cdot (y_i - \mathbf{x}_i^T \boldsymbol{\beta}_M) = 0$$

in which  $\boldsymbol{\beta}_M$  is the method of moments estimator of  $\boldsymbol{\beta}$ .

Since we have  $p$  predictors in total, and for each  $j = 1, \dots, p$  the same pattern follows, we can combine them into a system of  $p$  equations with number of unknowns being the dimension of  $\boldsymbol{\beta}_M$ , which is also  $p$ . We thus have the following equation:

$$\sum_{i=1}^n \mathbf{x}_i \cdot (y_i - \mathbf{x}_i^T \boldsymbol{\beta}_M) = 0 \Rightarrow \sum \mathbf{x}_i y_i - \sum \mathbf{x}_i \mathbf{x}_i^T \boldsymbol{\beta}_M = 0 \Rightarrow \sum \mathbf{x}_i y_i = \sum \mathbf{x}_i \mathbf{x}_i^T \boldsymbol{\beta}_M$$

Notice that the last equation is in exactly the same form as Equation (1) under **Least Squares** section. Therefore, the same conclusion follows, as  $\hat{\boldsymbol{\beta}} = \boldsymbol{\beta}_M$ . ■

(B)

*Sol:* We first derive the estimator with scalar  $w_i$ 's as weights as follows:

$$\begin{aligned}
\hat{\beta} &= \arg \min_{\beta} \left\{ \sum w_i (y_i - \mathbf{x}_i^T \beta)^2 \right\} \\
&\equiv \arg \min_{\beta} \begin{bmatrix} (y_1 - \mathbf{x}_1^T \beta)^T & (y_2 - \mathbf{x}_2^T \beta)^T & \dots & (y_n - \mathbf{x}_n^T \beta)^T \end{bmatrix} \begin{bmatrix} w_1 \\ \vdots \\ w_n \end{bmatrix} \begin{bmatrix} y_1 - \mathbf{x}_1^T \beta \\ y_2 - \mathbf{x}_2^T \beta \\ \vdots \\ y_n - \mathbf{x}_n^T \beta \end{bmatrix} \\
&\equiv \arg \min_{\beta} (\mathbf{Y} - \mathbf{X}\beta)^T \mathbf{W} (\mathbf{Y} - \mathbf{X}\beta) \quad (\text{define the diagonal matrix from above as } \mathbf{W}_{n \times n}) \\
&\equiv \arg \min_{\beta} (\mathbf{Y}^T \mathbf{W} \mathbf{Y} - 2\mathbf{Y}^T \mathbf{W} \mathbf{X} \beta + \beta^T \mathbf{X}^T \mathbf{W} \mathbf{X} \beta)
\end{aligned}$$

We apply similar arguments as those from **Least Squares** section below to derive  $\hat{\beta}$ :

The gradient of  $\beta$  from the last expression above is

$$\frac{\partial}{\partial \beta} = -2\mathbf{X}^T \mathbf{W} \mathbf{Y} + 2\mathbf{X}^T \mathbf{W} \mathbf{X} \beta$$

and the Hessian matrix is

$$\mathcal{H}(\beta) = 2\mathbf{X}^T \mathbf{W} \mathbf{X}$$

Given any non-zero vector  $\mathbf{u}$ , we have

$$\mathbf{u}^T \mathbf{X}^T \mathbf{W} \mathbf{X} \mathbf{u} = (\mathbf{X}\mathbf{u})^T \mathbf{W}^{1/2} \mathbf{W}^{1/2} (\mathbf{X}\mathbf{u}) = (\mathbf{W}^{1/2} \mathbf{X}\mathbf{u})^T (\mathbf{W}^{1/2} \mathbf{X}\mathbf{u})$$

As we pointed out before, since  $\mathbf{X}$  has full column rank,  $\mathbf{X}\mathbf{u}$  is a non-zero vector; and assuming none of the weights is 0, which is realistic (otherwise we don't even need this observation), we have non-zero vector  $\mathbf{W}^{1/2} \mathbf{X}\mathbf{u}$ . Therefore,

$$\mathbf{u}^T \mathbf{X}^T \mathbf{W} \mathbf{X} \mathbf{u} = (\mathbf{W}^{1/2} \mathbf{X}\mathbf{u})^T (\mathbf{W}^{1/2} \mathbf{X}\mathbf{u}) > 0$$

i.e.  $\mathbf{X}^T \mathbf{W} \mathbf{X}$  is positive definite.

Therefore, by solving for  $\beta$  when gradient equals to 0 we have the critical point of the original summation to be minimized at  $(\mathbf{X}^T \mathbf{W} \mathbf{X})^{-1} \mathbf{X}^T \mathbf{W} \mathbf{Y}$ , and it gives the minimum of this summation due to the positive-definiteness of the Hessian matrix, i.e.

$$\hat{\beta} = (\mathbf{X}^T \mathbf{W} \mathbf{X})^{-1} \mathbf{X}^T \mathbf{W} \mathbf{Y}$$

We now derive the estimator with heteroscedastic Gaussian error:

$$\begin{aligned}
\hat{\beta} &= \arg \max_{\beta} \left\{ \prod p(y_i) \right\} \\
&\equiv \arg \max_{\beta} \left\{ \log \prod p(y_i) \right\} \\
&\equiv \arg \min_{\beta} \left\{ \sum \frac{1}{\sigma_i^2} (y_i - \mathbf{x}_i^T \beta)^2 \right\}
\end{aligned}$$

We notice from the last expression s.t. by setting each

$$\frac{1}{\sigma_i^2} = w_i$$

, the exact same arguments as our derivation above would follow. ■

# Quantifying Uncertainty: Some Basic Frequentist Ideas

## In linear Regression

(A)

*Sol:* We again use mgf to derive the sampling distribution of  $\hat{\beta}$  (not sure about the term "sampling"; here I'm deriving "the" distribution of  $\hat{\beta}$ ):

$$\begin{aligned}
M_{\hat{\beta}}(\mathbf{t}) &= \mathbb{E}[e^{\mathbf{t}^T \hat{\beta}}] \\
&= \mathbb{E}[e^{\mathbf{t}^T ((\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{Y})}] \\
&= \mathbb{E}[e^{\mathbf{t}^T ((\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T (\mathbf{X}\beta + \epsilon))}] \\
&= \mathbb{E}[e^{\mathbf{t}^T ((\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{X}\beta + (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \epsilon)}] \\
&= e^{\mathbf{t}^T \beta} \cdot \mathbb{E}[e^{\mathbf{t}^T ((\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \epsilon)}] \\
&= e^{\mathbf{t}^T \beta} \cdot M_\epsilon(\mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{t}) \\
&= \exp(\mathbf{t}^T \beta) \cdot \exp((0 + (\mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{t})^T (\mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{t}) \sigma^2 / 2)) \\
&= \exp(\mathbf{t}^T \beta + \mathbf{t}^T (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{t} \sigma^2 / 2)
\end{aligned}$$

Note that the last expression is the mgf of a  $\mathcal{MVN}(\beta, (\mathbf{X}^T \mathbf{X})^{-1} \sigma^2)$  random variable. Since a mgf uniquely determines the distribution, we have

$$\hat{\beta} \sim \mathcal{MVN}(\beta, (\mathbf{X}^T \mathbf{X})^{-1} \sigma^2)$$

■

(B)

*Sol:* To quantify the uncertainty of each of the  $p$  coefficients (including one for the intercept here), ideally we can simply read off the diagonal of  $\hat{\beta}$ 's covariance matrix. However, this covariance matrix includes the term  $\sigma^2$ , which in our case we don't know its value. Therefore, we need to have an estimator of  $\sigma^2$  from our sample data, and it would be nice if such estimator is unbiased.

Fortunately, we've been shown from last semester's Linear Models class such that, denoting residuals as  $\hat{\epsilon}$  such that

$$\hat{\epsilon} = \mathbf{Y} - \hat{\mathbf{Y}}$$

we'd have

$$\hat{\epsilon}^T \hat{\epsilon} \sim \sigma^2 \chi^2_{n-p}$$

Therefore, its expectation is

$$\mathbb{E}(\hat{\epsilon}^T \hat{\epsilon}) = \sigma^2 \mathbb{E}(\text{a } \chi^2_{n-p} \text{ random variable}) = \sigma^2 \cdot (n - p)$$

We thus have an unbiased estimator of  $\sigma^2$ :

$$s^2 = \frac{\hat{\epsilon}^T \hat{\epsilon}}{n - p}$$

We thus can estimate the covariance matrix of  $\hat{\beta}$  by

$$S^2 = (\mathbf{X}^T \mathbf{X})^{-1} \cdot s^2 = (\mathbf{X}^T \mathbf{X})^{-1} \frac{\hat{\epsilon}^T \hat{\epsilon}}{n-p} = (\mathbf{X}^T \mathbf{X})^{-1} \frac{(\mathbf{Y} - \hat{\mathbf{Y}})^T (\mathbf{Y} - \hat{\mathbf{Y}})}{n-p}$$

We've implemented such computation in R, and here's the code snippet:

```
# Fill in the blank
# compute residual
residual = y - x %*% betahat
res_T_res = t(residual) %*% residual
# calculate unbiased estimator of sigma squared
s_squared = (res_T_res/(dim(x)[1]-dim(x)[2]))[1,1]
# estimate covariance matrix of beta
betacov = solve(t(x) %*% x) * s_squared
# calculate standard error by hand
beta_se = sqrt(diag(betacov))
```

We found that our calculation of standard errors for each  $\beta_j$  exactly matches those from lm library:

```
> beta_se = sqrt(diag(betacov))
> beta_se_lm = sqrt(diag(betacov_lm))
> beta_se
      V5          V6          V7          V8          V9          V10
38.3286939914 0.0072511036 0.1741433372 0.0237693597 0.0692994207 0.1247136189 0.0003943697
      V11         V12         V13
0.0147771925 0.1192916644 0.0048962536
> beta_se_lm
      x          xV5          xV6          xV7          xV8          xV9          xV10
38.3286939914 0.0072511036 0.1741433372 0.0237693597 0.0692994207 0.1247136189 0.0003943697
      xV11        xV12        xV13
0.0147771925 0.1192916644 0.0048962536
> |
```

My guess is that if this matches after this many decimal spaces, then lm's implementation might be exactly the same as the one we provided. ■

## Propagating Uncertainty

(A)

*Sol:* Denoting the standard error of  $f(\hat{\theta})$  as  $se(f(\hat{\theta}))$ , we have

$$\begin{aligned} [se(f(\hat{\theta}))]^2 &= var(f(\hat{\theta})) \\ &= \mathbb{E}[(f(\hat{\theta}) - \mathbb{E}f(\hat{\theta}))^2] \\ &= \mathbb{E}[((\hat{\theta}_1 + \hat{\theta}_2) - (\theta_1 + \theta_2))^2] \\ &= \mathbb{E}[((\hat{\theta}_1 - \theta_1) + (\hat{\theta}_2 - \theta_2))^2] \\ &= \mathbb{E}[(\hat{\theta}_1 - \theta_1)^2 + (\hat{\theta}_2 - \theta_2)^2 + 2(\hat{\theta}_1 - \theta_1)(\hat{\theta}_2 - \theta_2)] \\ &= var(\hat{\theta}_1) + var(\hat{\theta}_2) + 2 \cdot cov(\hat{\theta}_1, \hat{\theta}_2) \\ &= \hat{\Sigma}_{11} + \hat{\Sigma}_{22} + 2 \cdot \hat{\Sigma}_{12} \end{aligned}$$

thus the standard error is just taking the square root of our estimated variance of such linear function of estimated parameter.

We directly write out the square of standard error for  $f(\hat{\theta}) = \sum_{j=1}^p \hat{\theta}_j$ :

$$[se(f(\hat{\theta}))]^2 = var(f(\hat{\theta})) = \sum_{j=1}^p \hat{\Sigma}_{jj} + \sum_{1 \leq i < j \leq p} 2 \cdot \hat{\Sigma}_{ij}$$

after which we take the square root to get the standard error.

Using matrix notation: define

$$\mathbf{J}_{p \times 1} = [1 \quad 1 \quad \dots \quad 1]^T$$

we have

$$var(f(\hat{\theta})) = var\left(\sum_{j=1}^p \hat{\theta}_j\right) = var(\mathbf{J}^T \hat{\theta}) = \mathbf{J}^T var(\hat{\theta}) \mathbf{J}$$

in which the  $(i, j)$ -th element in  $var(\hat{\theta})$  is  $\hat{\Sigma}_{ij}$ . ■

## (B)

*Sol:* The first order Taylor approximation of  $f(\hat{\theta})$  around the unknown true parameter value  $\theta$  is

$$f(\hat{\theta}) = f(\theta) + \nabla_{\hat{\theta}}(\theta)(\hat{\theta} - \theta)$$

Therefore, we approximate the variance of  $f(\hat{\theta})$  by

$$\begin{aligned} var(f(\hat{\theta})) &= var(f(\theta) + \nabla_{\hat{\theta}}(\theta)(\hat{\theta} - \theta)) \\ &= var(\nabla_{\hat{\theta}}(\theta)(\hat{\theta} - \theta)) \\ &= \nabla_{\hat{\theta}}(\theta) var(\hat{\theta} - \theta) \nabla_{\hat{\theta}}^T(\theta) \\ &= \nabla_{\hat{\theta}}(\theta) var(\hat{\theta}) \nabla_{\hat{\theta}}^T(\theta) \\ &= \nabla_{\hat{\theta}}(\theta) \hat{\Sigma} \nabla_{\hat{\theta}}^T(\theta) \end{aligned}$$

in which the simplification of above steps uses the fact that true parameter  $\theta$  can be treated as a constant in the function of  $\hat{\theta}$ .

Since the true parameter is unknown, we are not able to calculate the value  $\nabla_{\hat{\theta}}(\theta)$ , therefore we use our estimated value  $\hat{\theta}$  in the place of  $\theta$  and have

$$var(f(\hat{\theta})) = \nabla_{\hat{\theta}}(\hat{\theta}) \hat{\Sigma} \nabla_{\hat{\theta}}^T(\hat{\theta})$$

We observe three estimations:

- Taylor approximation of  $f$
- approximation of gradient at unknown true  $\theta$  using the value of  $\hat{\theta}$
- approximation of covariance matrix by using square of standard error in place of  $\sigma^2$

Each estimation could add extra deviation from true  $var(f(\hat{\theta}))$ . ■